

Guest Lecture for 6.869 Advances in Computer Vision

# Activity Recognition

Bolei Zhou

MIT CSAIL

# Challenge for Image Recognition

- Variation in appearance.



# Challenge for Activity Recognition

- Describing activity at the proper level

Image recognition?  
No motion needed?

Skeleton recognition?  
Which activities?



# Challenge for Activity Recognition

- Describing activity at the proper level

A chain of events

Making chocolate cookies



# Outline

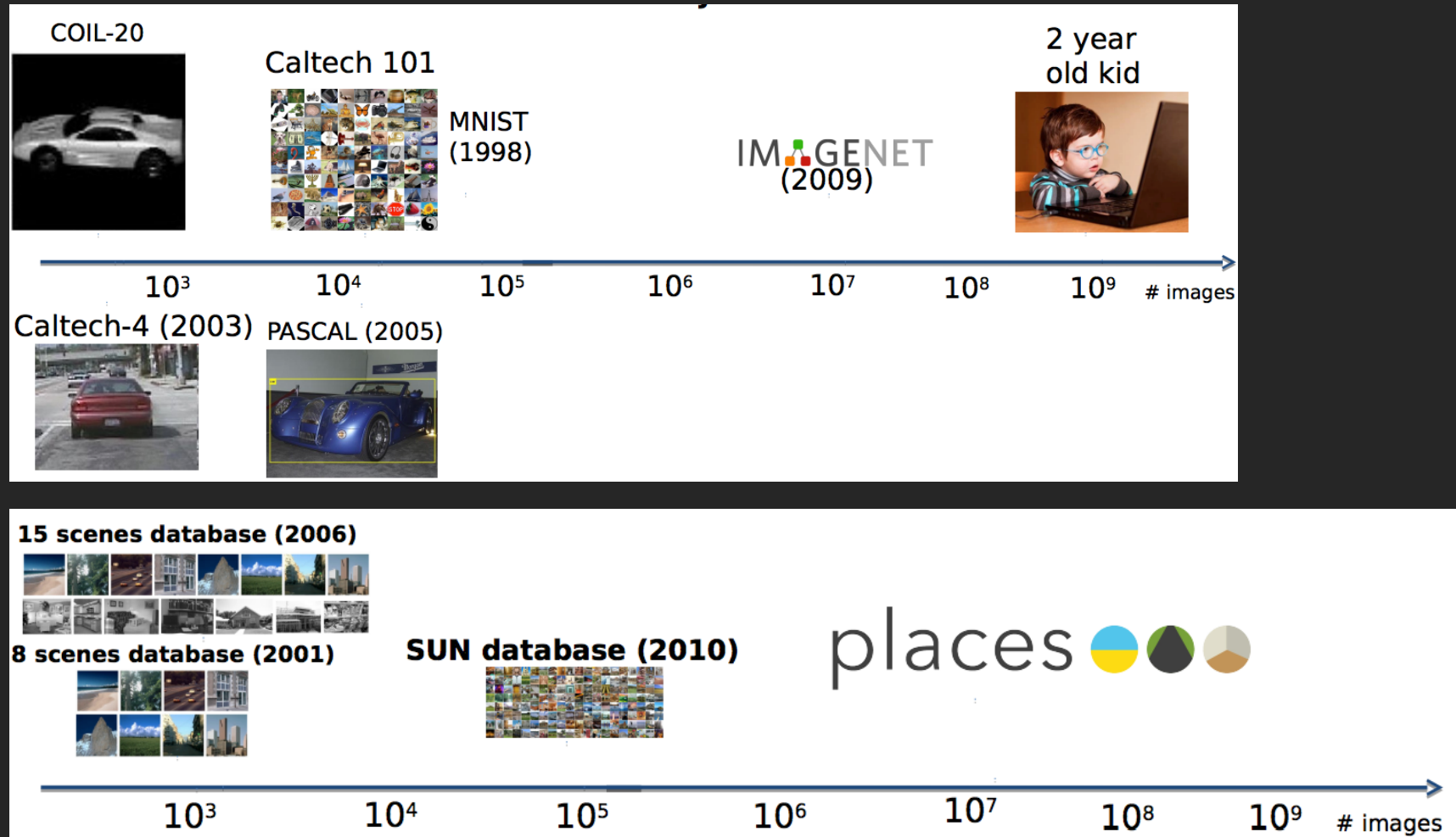
- Video Recognition Datasets
- Video Recognition Models

A little bit about my recent work:

- Temporal Relational Reasoning in Videos

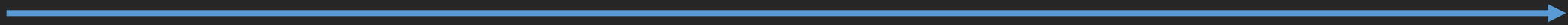
# Video Recognition Datasets

- Review on image datasets



# Video Recognition Datasets

KTH:6      HMDB:51      Kinetics:400  
UCF: 101      Moments: 339  
ActivityNet:200



# Video Recognition Datasets

Two video collection methods:

- Collect videos from the web (Youtube, Flickr, etc)
- Crowd-sourcing video collection.



# Video Recognition Datasets

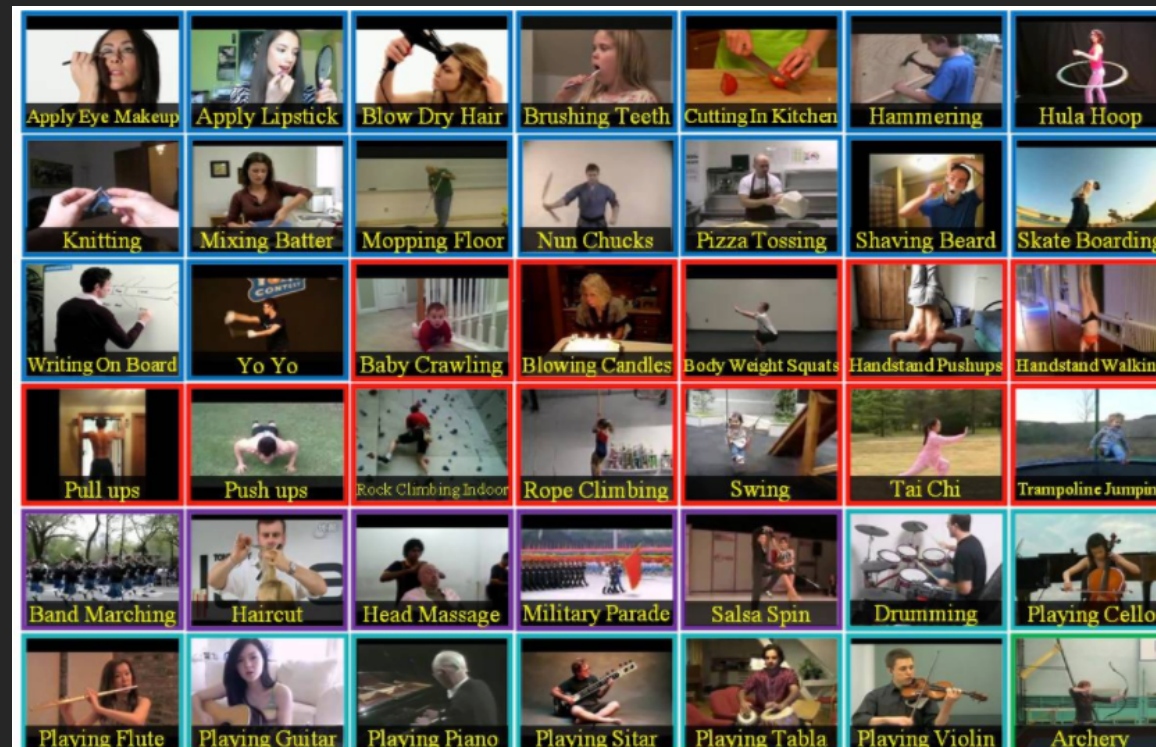
- KTH Dataset: recognition of human actions
- 6 classes, 2391 videos <https://www.youtube.com/watch?v=Jm69kbCC17s>



# Video Recognition Datasets

- UCF101 from University of Central Florida
- 101 classes, 9,511 videos in training

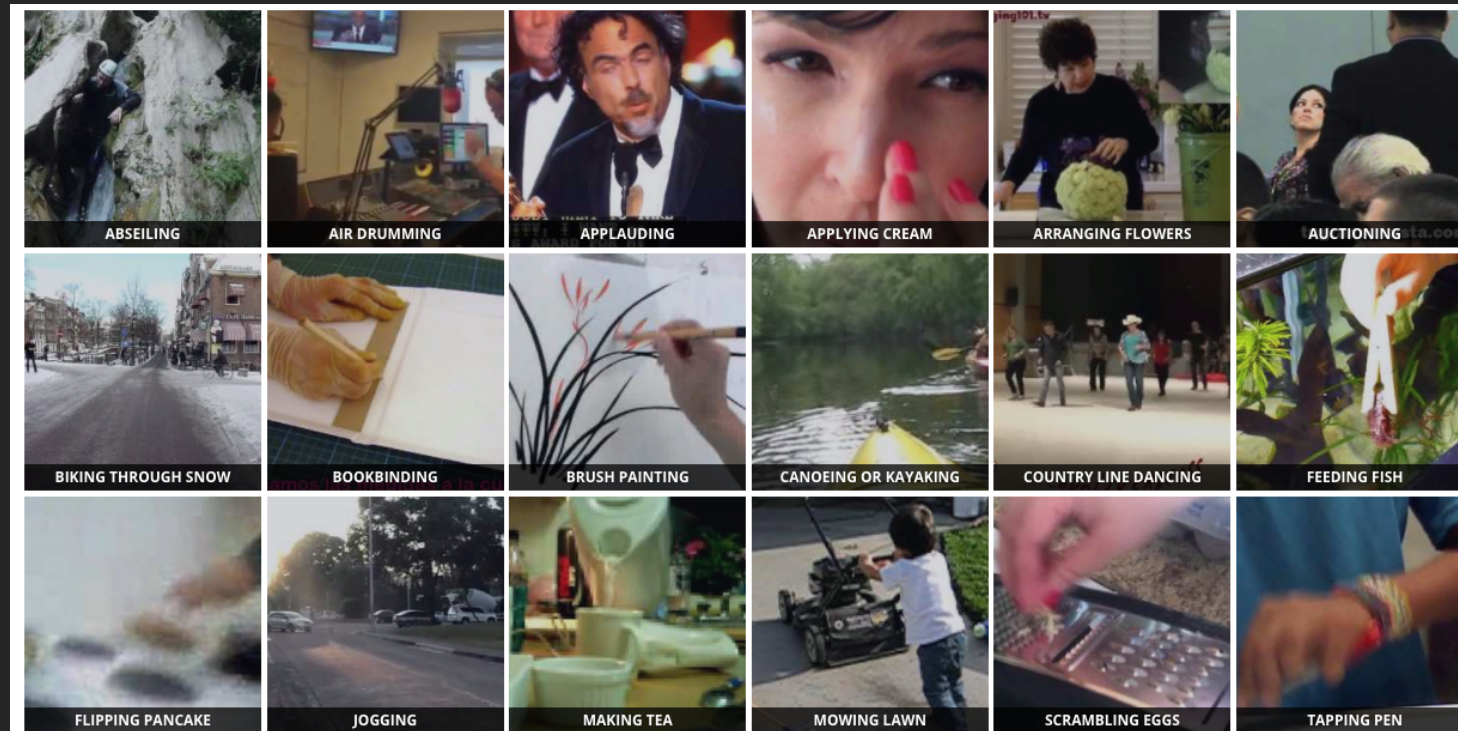
<https://www.youtube.com/watch?v=hGhuUaxocIE>



# Video Recognition Datasets

- Kinetics from Google DeepMind
- 400 classes, 239,956 videos in training

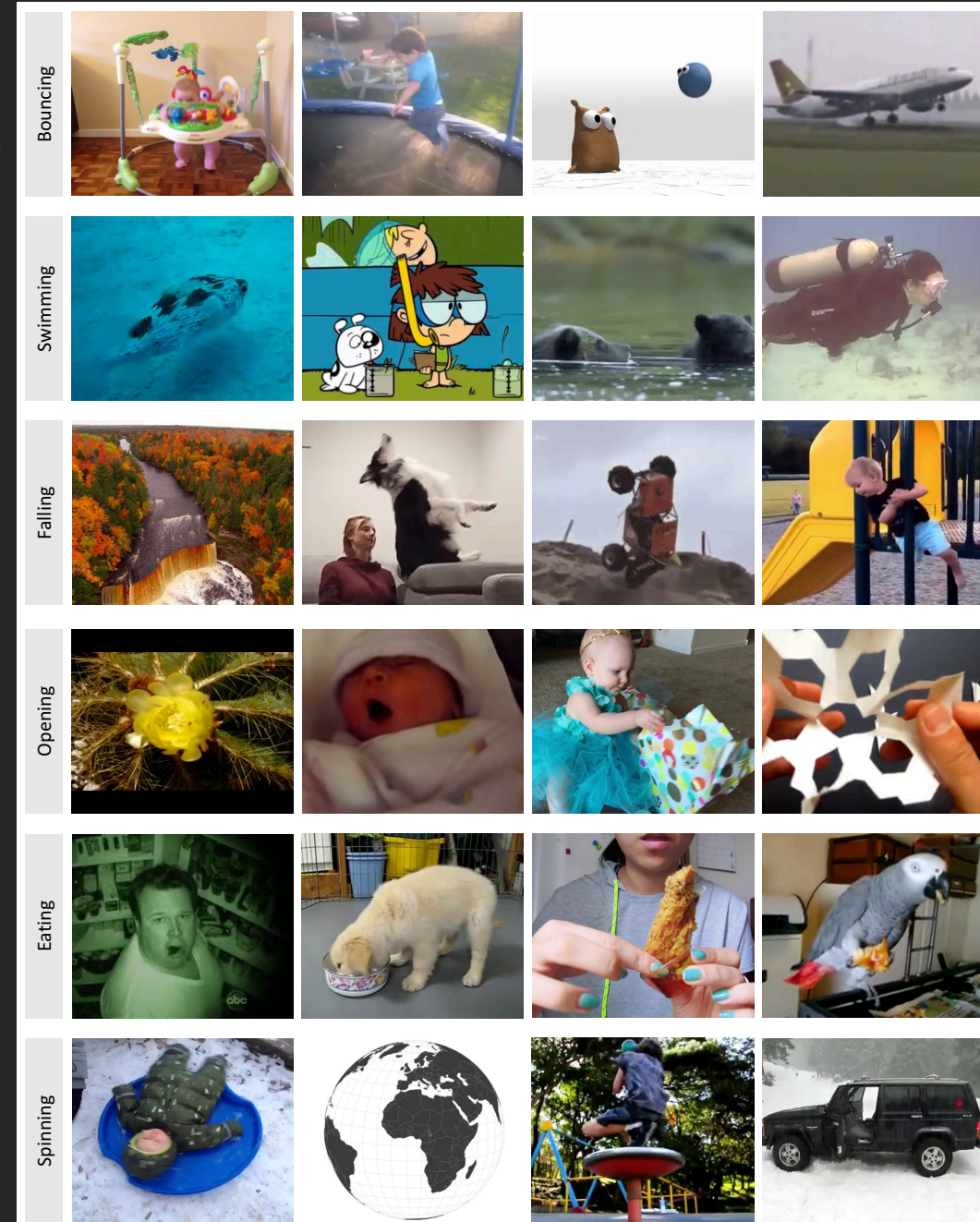
<https://deepmind.com/research/open-source/open-source-datasets/kinetics/>



# Video Recognition Datasets

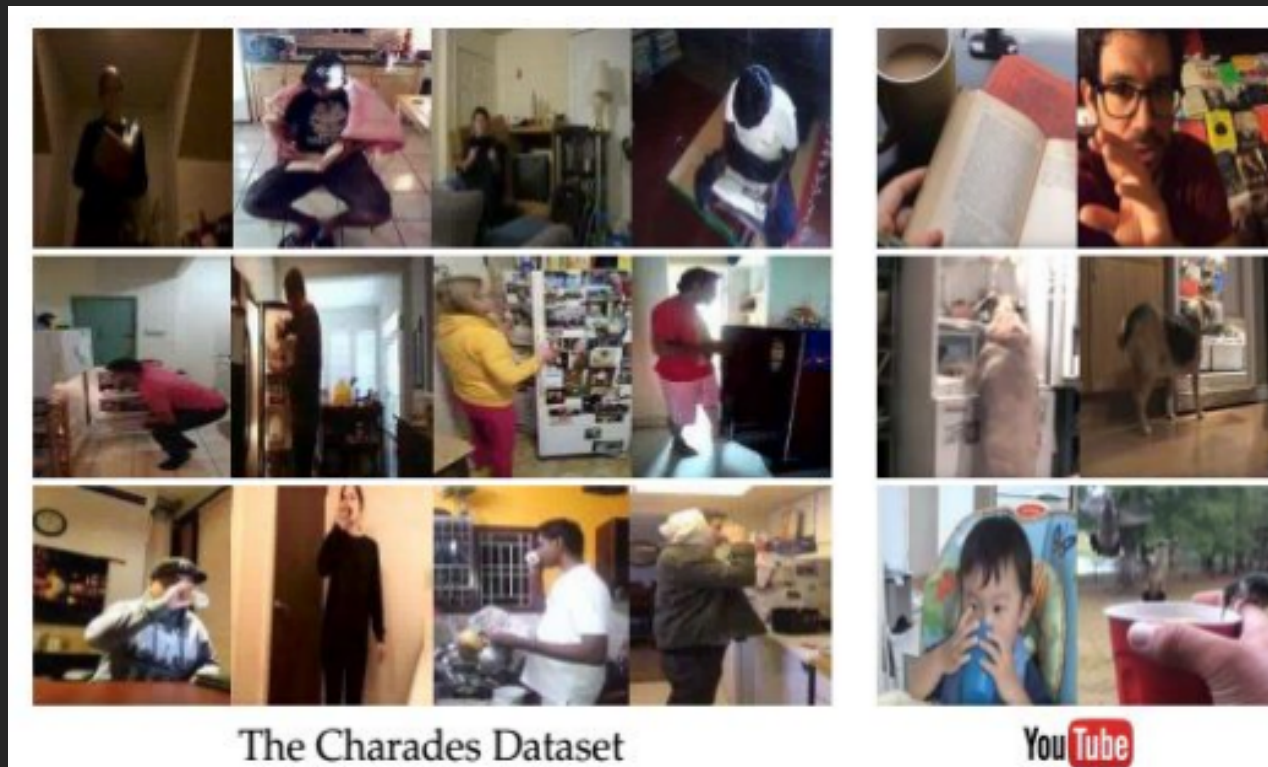
- Moments from MIT
- 1 million 3-second video from 339 generic actions

[http://moments.csail.mit.edu/index\\_test.html](http://moments.csail.mit.edu/index_test.html)



# Video Recognition Datasets

- Charades dataset: Hollywood in Homes
- Crowdsourced video dataset



# Video Recognition Datasets

- Charades dataset: Hollywood in Homes
- Long chain of actions

1st Generation



2nd Generation



3rd Generation



Now

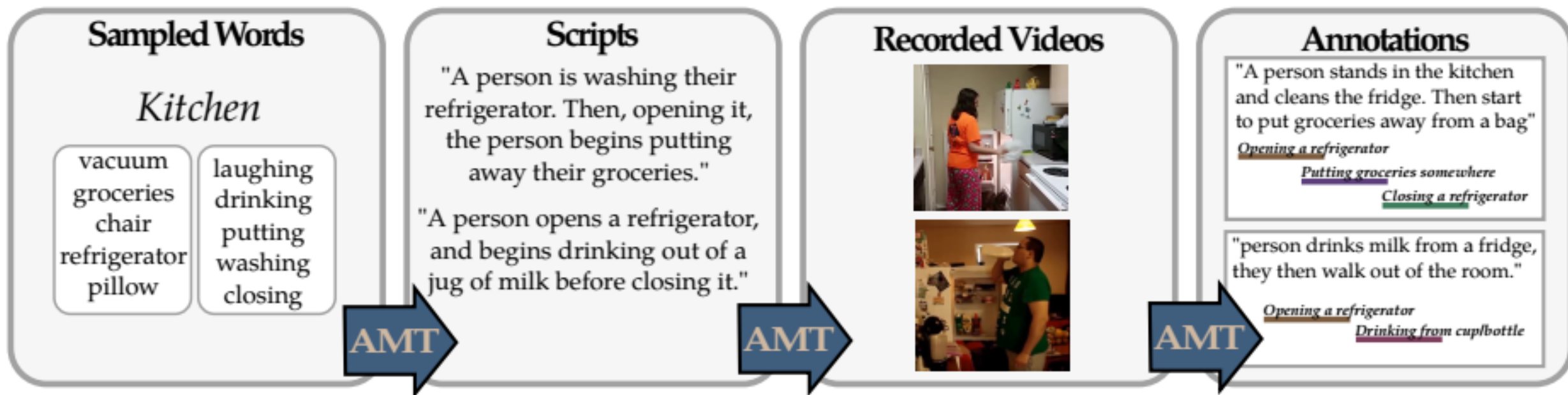


<https://www.youtube.com/watch?v=x9AhZLDkbyc>

# Video Recognition Datasets



- Charades dataset: Hollywood in Homes
- Crowd-sourced video dataset



# Video Recognition Datasets



- Charades dataset: Hollywood in Homes
- Demo video

<https://www.youtube.com/watch?v=x9AhZLDkbyc>



# Video Recognition Datasets

- Something-Something dataset: human object interaction
  - 174 categories: 100,000 videos
- 
- Holding something
  - Turning something upside down
  - Turning the camera left while filming something
  - Opening something



Poking a stack of something  
so the stack collapses



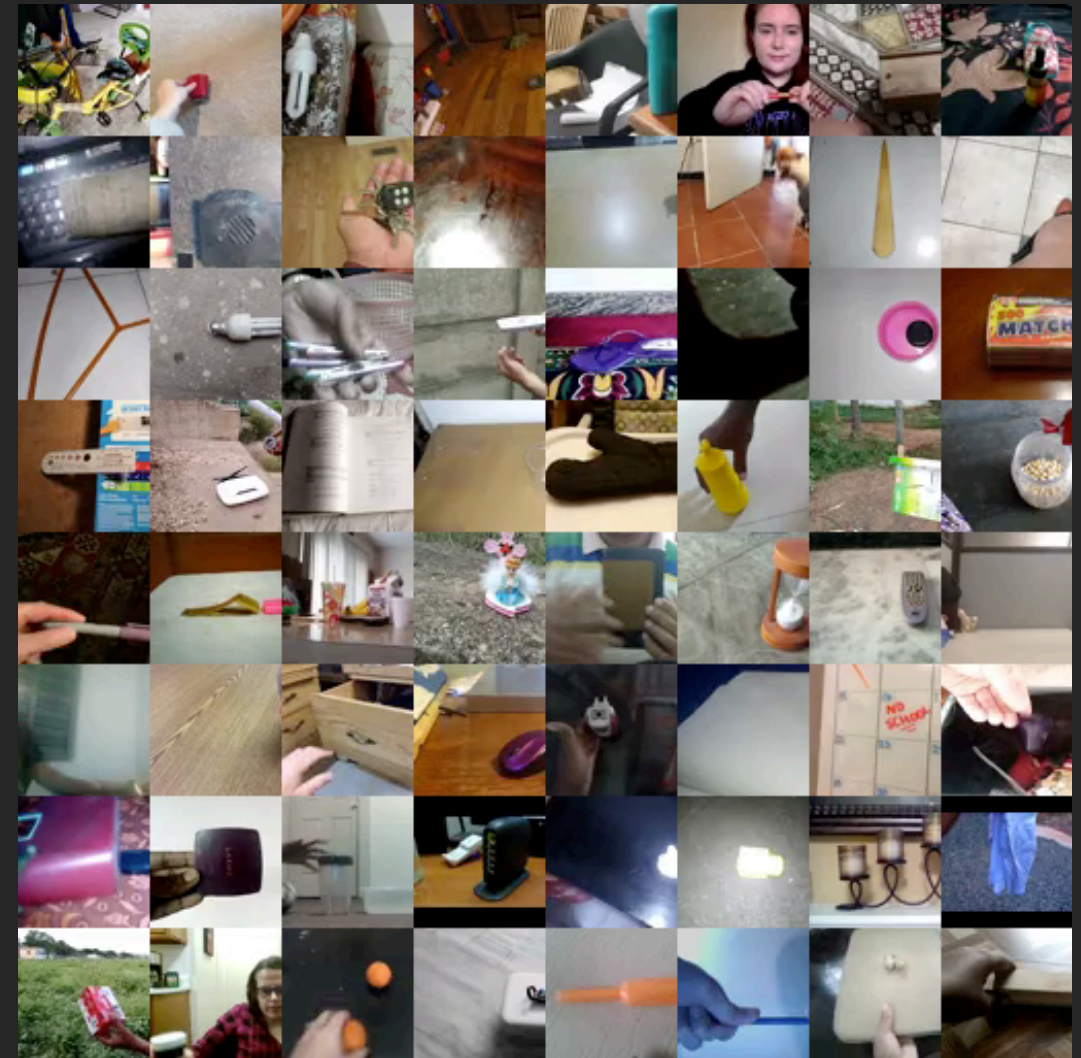
Plugging something into  
something

# Crowd-sourcing Video Collection



You have uploaded **0** of **10** videos **Submit Task**

|  |  |
|--|--|
| <p>holding something over something</p> <p><b>Upload Video</b></p>   | <p>moving something away from something</p> <p><b>Upload Video</b></p> |
| <p>lifting a surface with something on it but not enough for it to slide down</p> <p><b>Upload Video</b></p> | <p>moving something away from something</p> <p><b>Upload Video</b></p> |
| <p>tearing something into two pieces</p> <p><b>Upload Video</b></p>  |  |

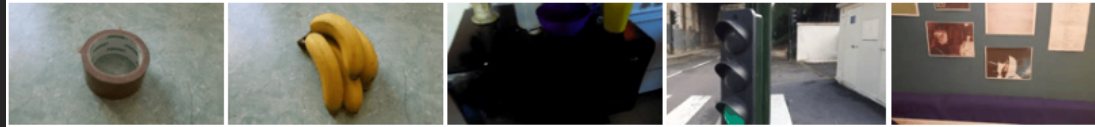


<https://www.twentybn.com/datasets/something-something>

# Something-to-Something

[http://visiongpu23.csail.mit.edu/deepscene/moments/models/datasets/something/plot\\_gif.html](http://visiongpu23.csail.mit.edu/deepscene/moments/models/datasets/something/plot_gif.html)

## 0 Approaching something with your camera



## 1 Attaching something to something



## 2 Bending something so that it deforms



## 3 Bending something until it breaks



# Video = Sequence of RGB images

How to represent temporal information?

- Capture the temporal dependency
- Efficiency: 1min 25fps video = 1500 images



# Video Recognition Models

- Pre-Deep learning era

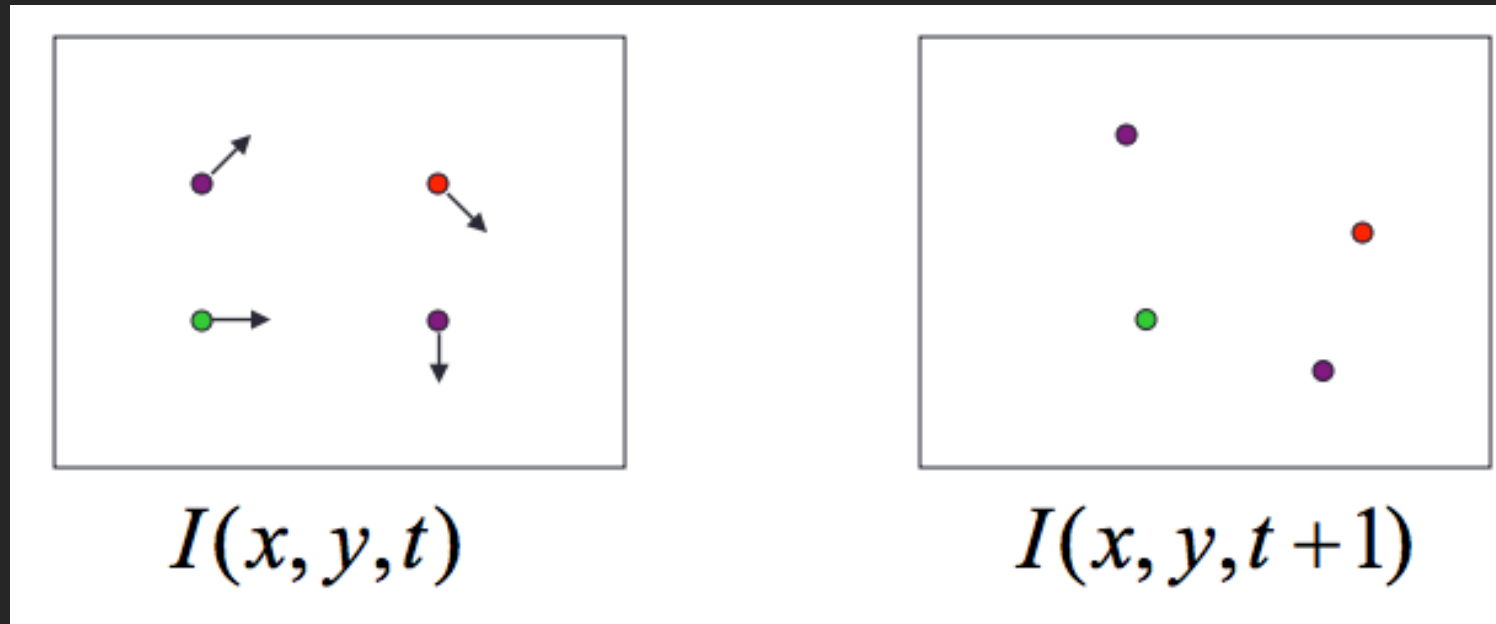
Optic flow, trajectories, bag of words.

- Deep learning era

Neural Networks

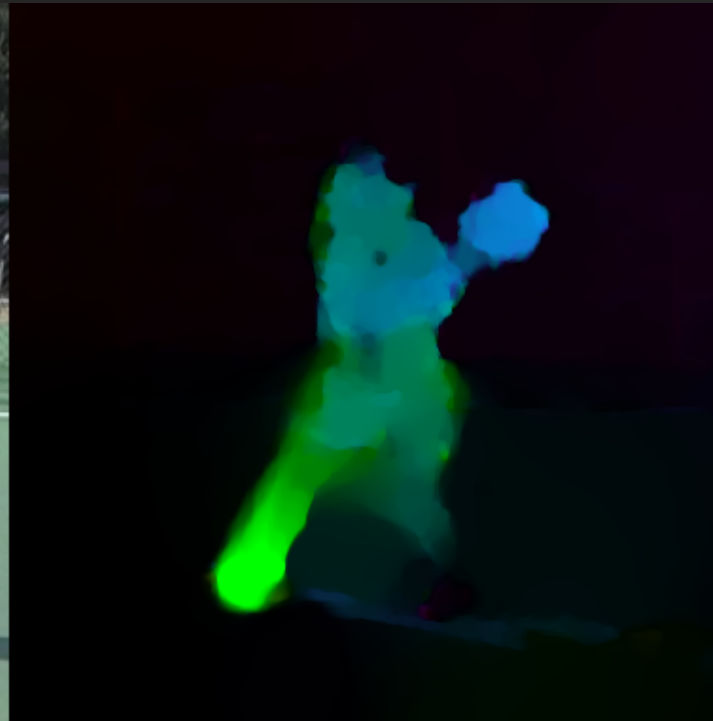
# Pre-deep learning Activity Recognition

- Optic Flow: the displacement of pixels
- Gesture lecture by Ce Liu next week on motion estimation



# Motion Representations in Activity Recognition

- Optic Flow



<https://www.youtube.com/watch?v=JSzUdVBmQP4>

# Motion Representations in Activity Recognition

- Trajectories: key-point tracking over frames



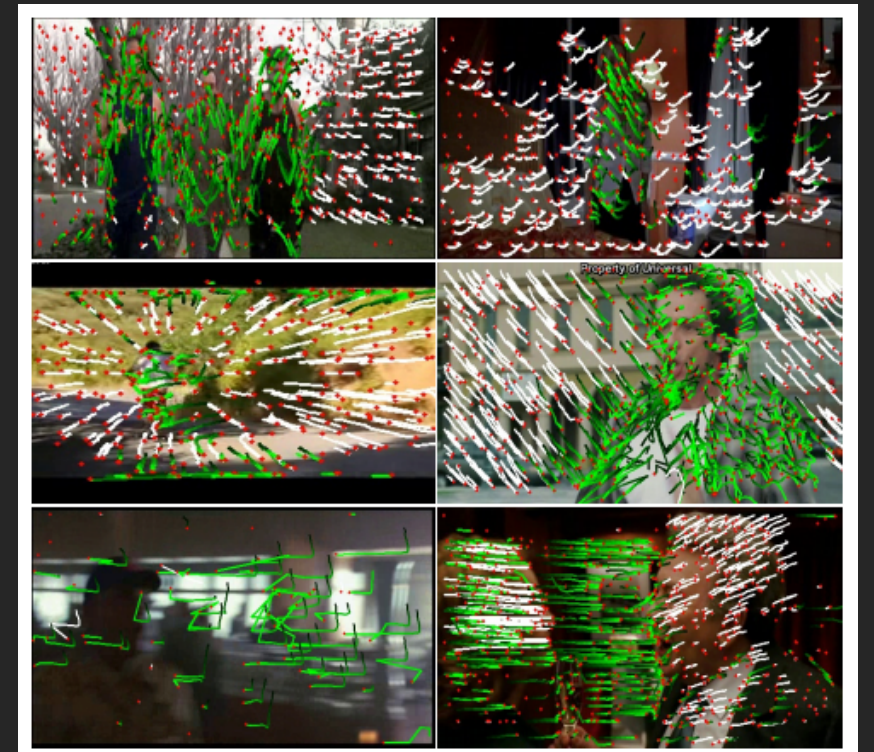
<https://www.youtube.com/watch?v=YN2lDqXz-uc>



# Motion Representations in Activity Recognition

## Improved Dense Trajectory (iDT)

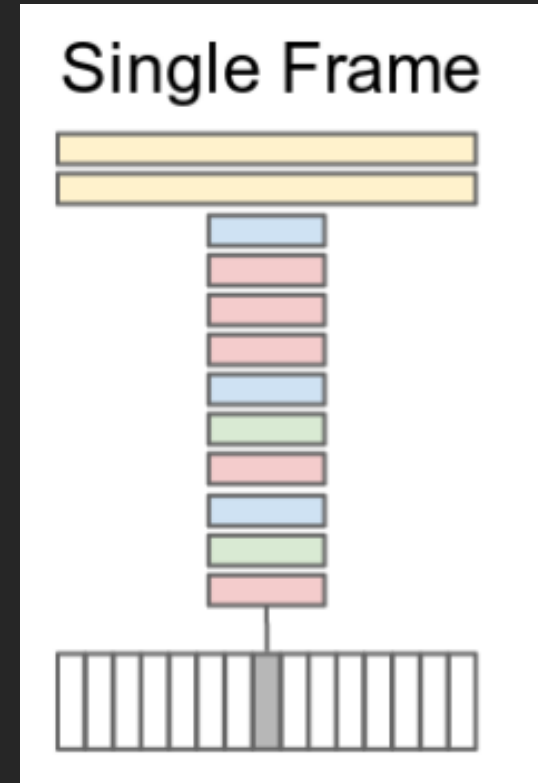
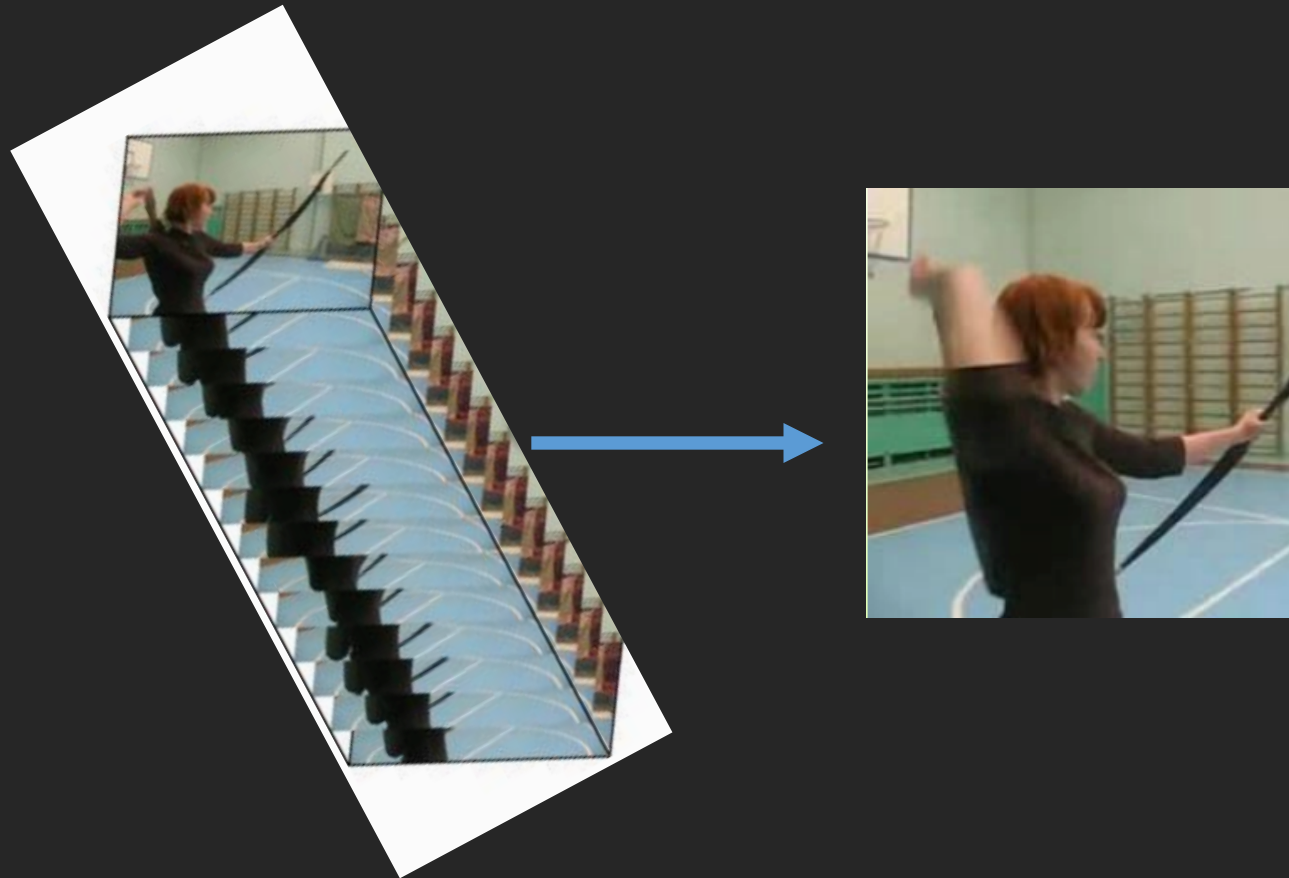
- Global motion compensation (camera motion removal)
- Features from trajectories and HoG
- Bag of trajectories + Fisher Vector + PCA



# Deep Learning Models for Activity Recognition

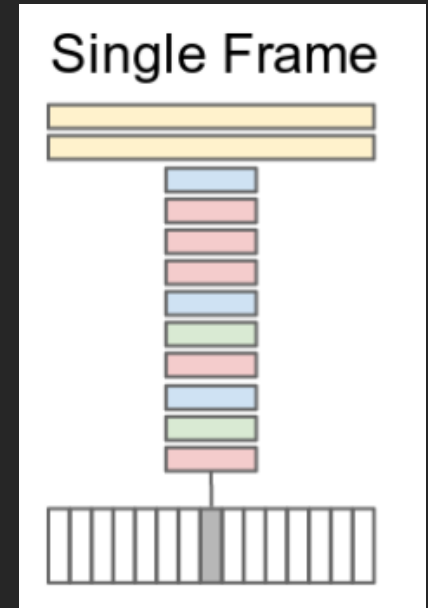
- RGB frame fusion network
- 2-stream network
- 3D convolution network
- Temporal segment network

# Single-frame image model

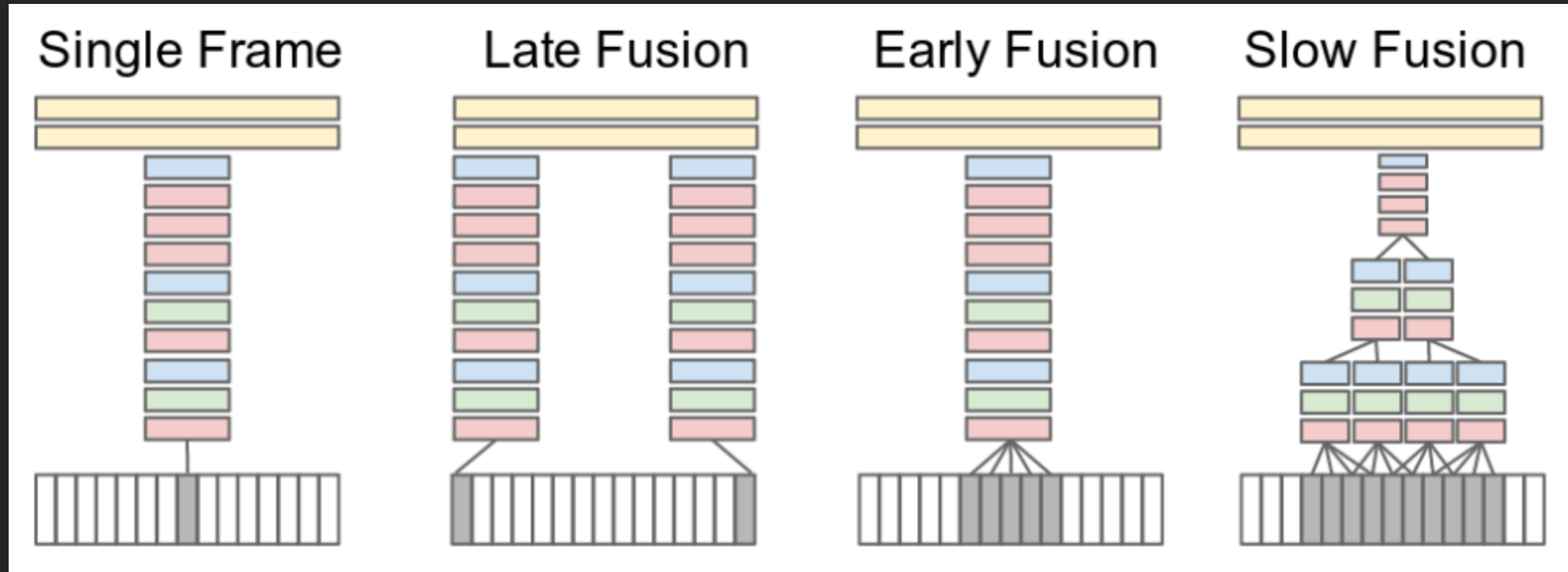


# Performance on the UCF101

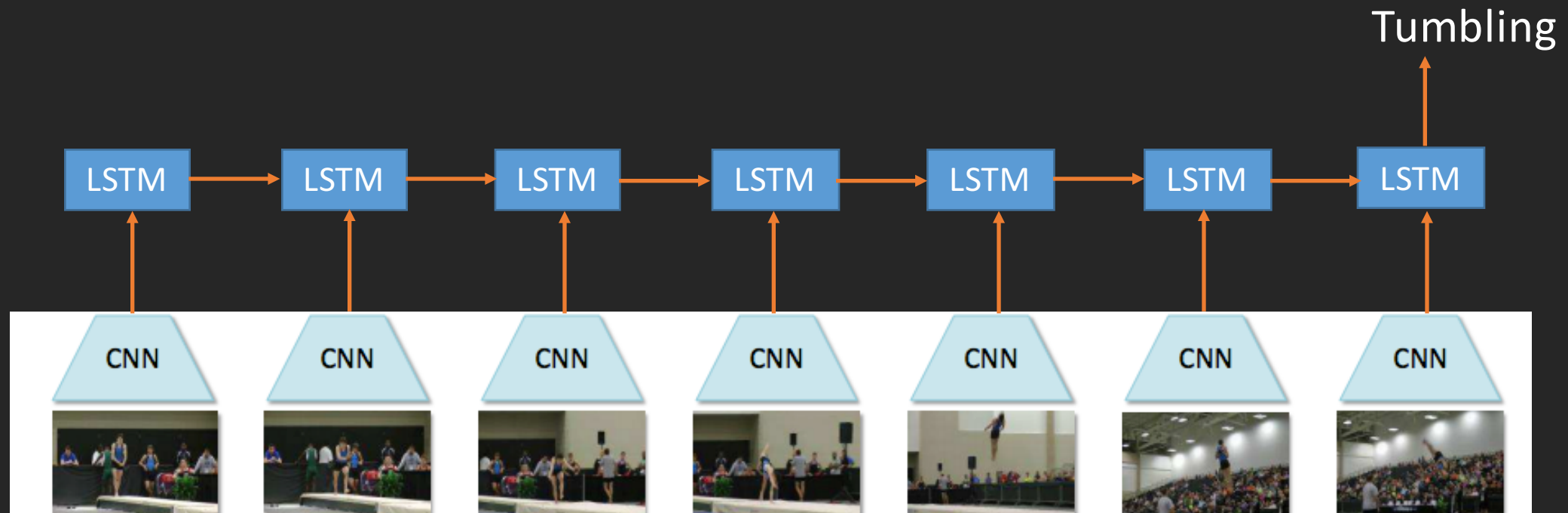
| Spatial ConvNets | Temporal ConvNets | Two-Stream |
|------------------|-------------------|------------|
| 72.7%            | 81.0%             | 87.0%      |



# Multi-frame fusion model

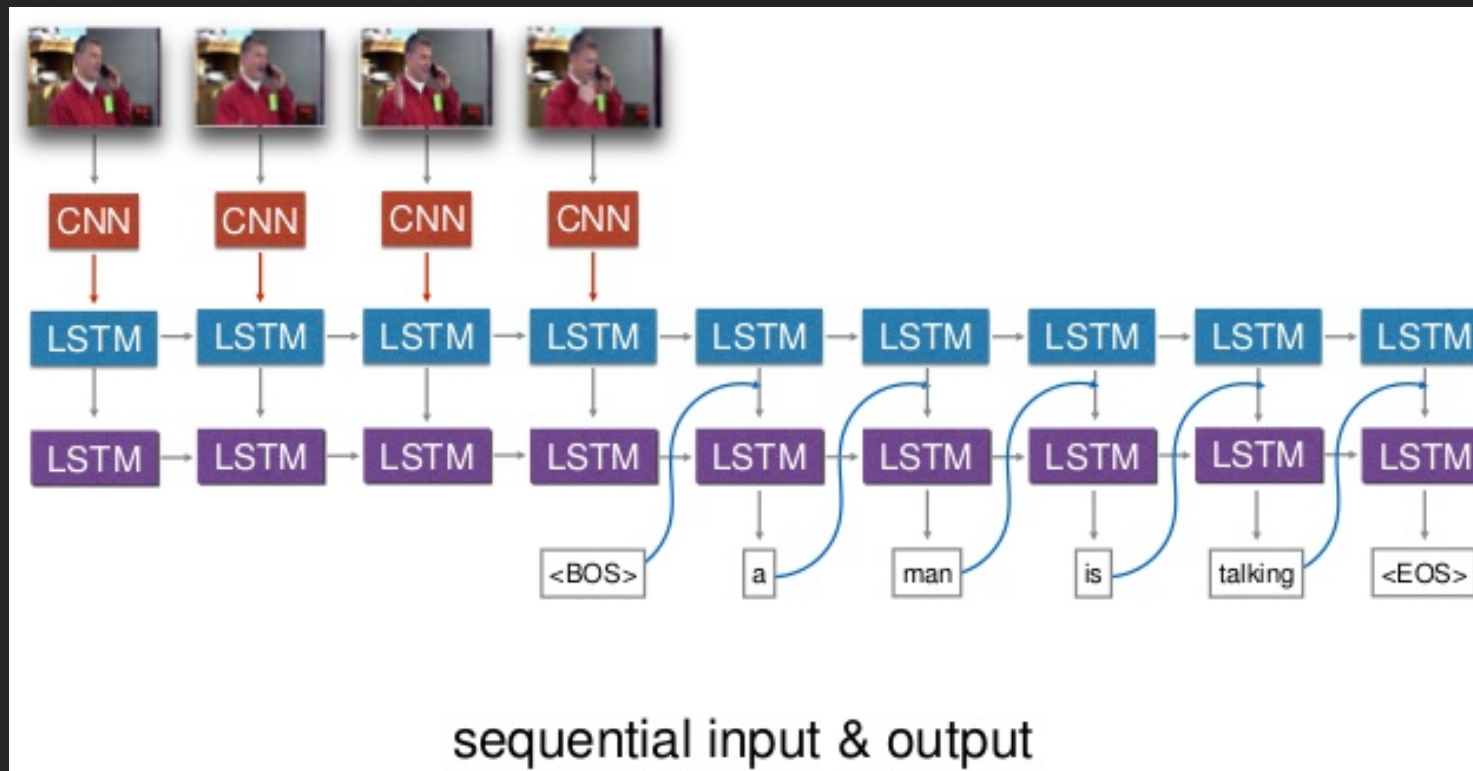


# Multi-frame LSTM fusion model

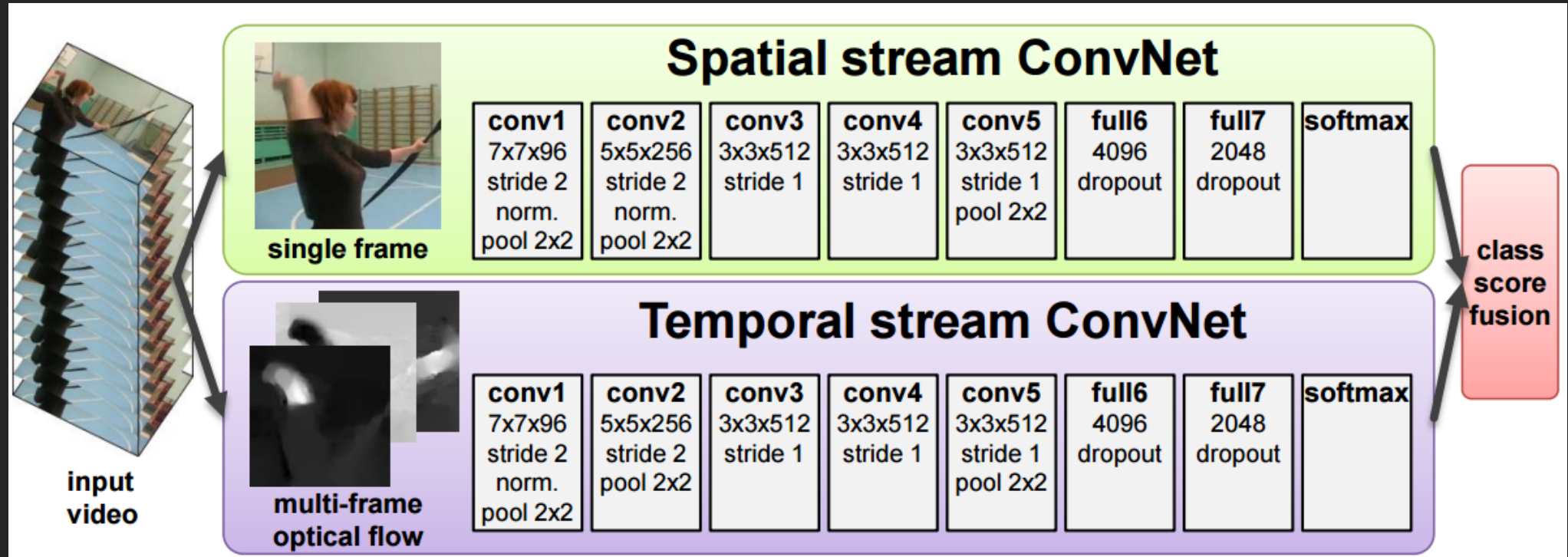


# LSTM: recursive neural networks

- Video Captioning

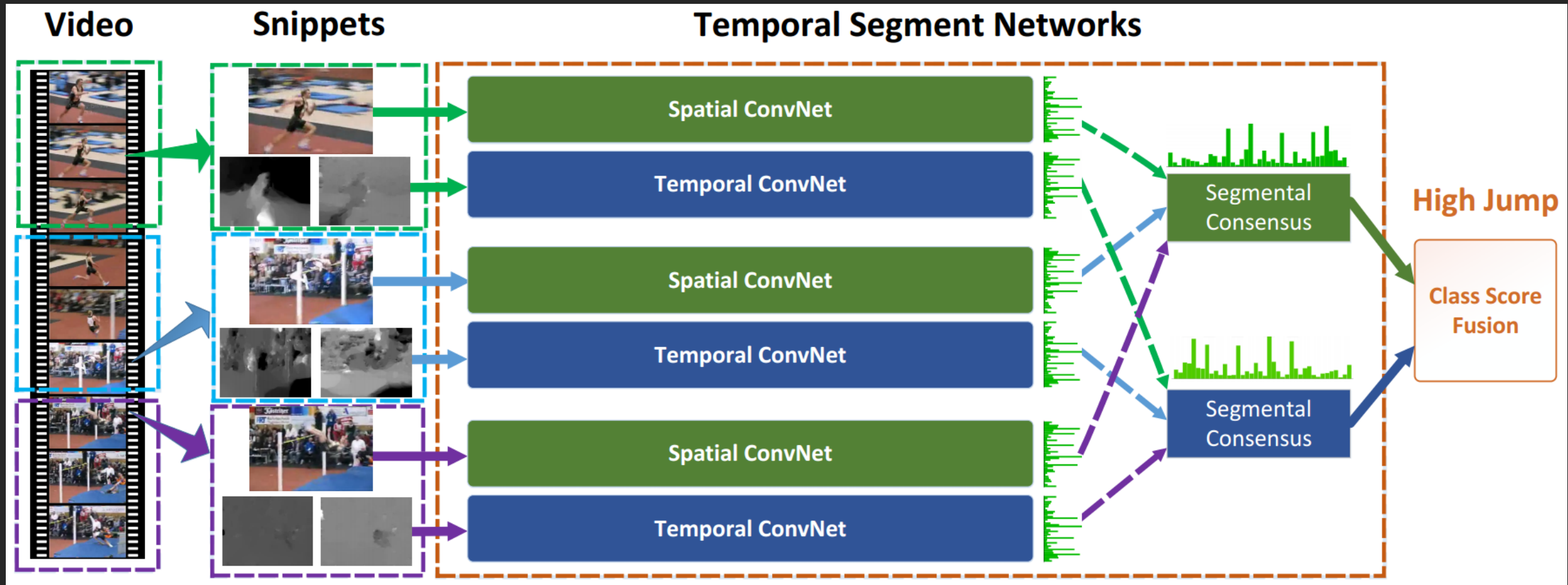


# 2-Stream Network



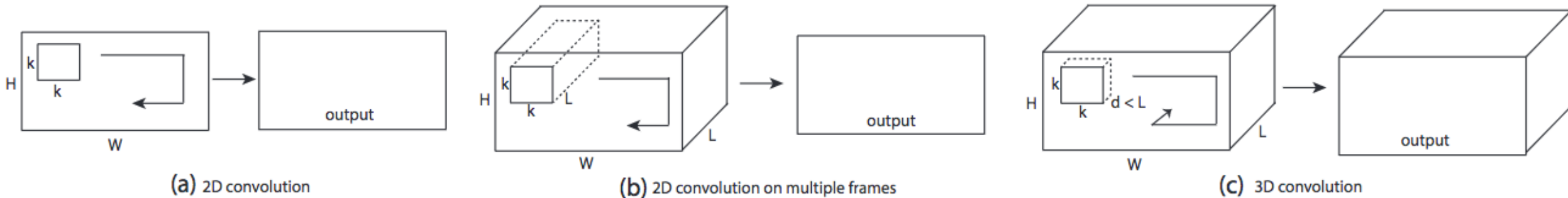


# Temporal segment network



# 3D convolutional Networks

Computationally expensive, and a lot of model parameters



If it is RGB frame rather than grey frame, it is actually 4D convolution.

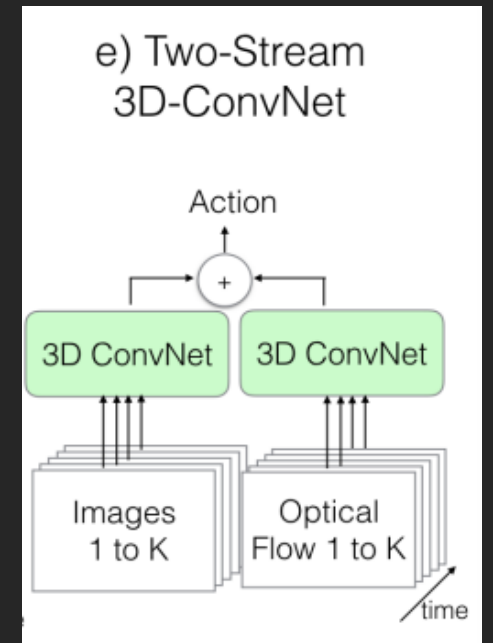
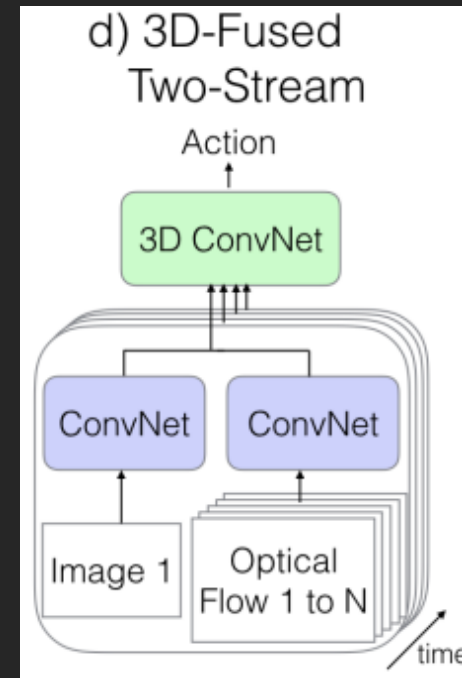
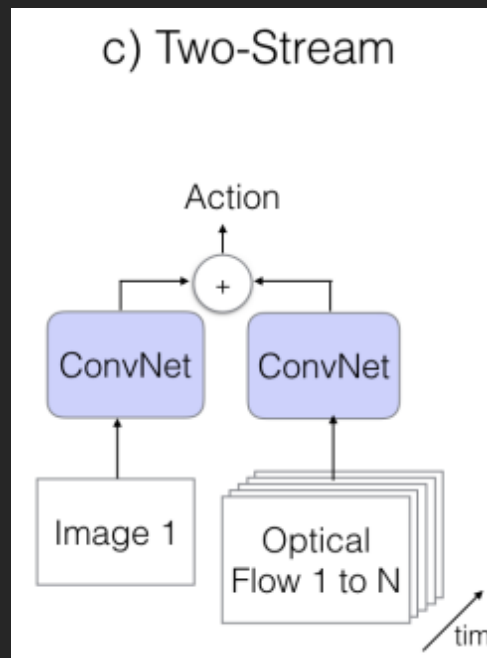
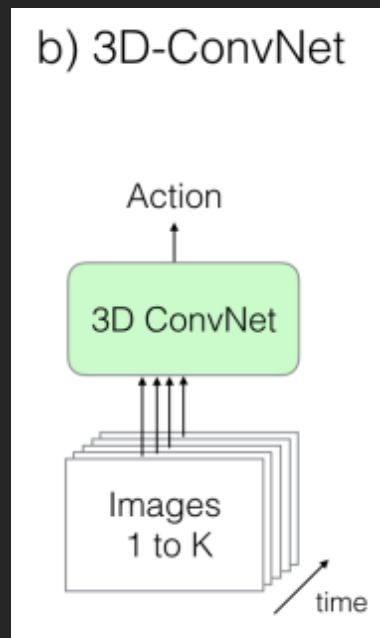
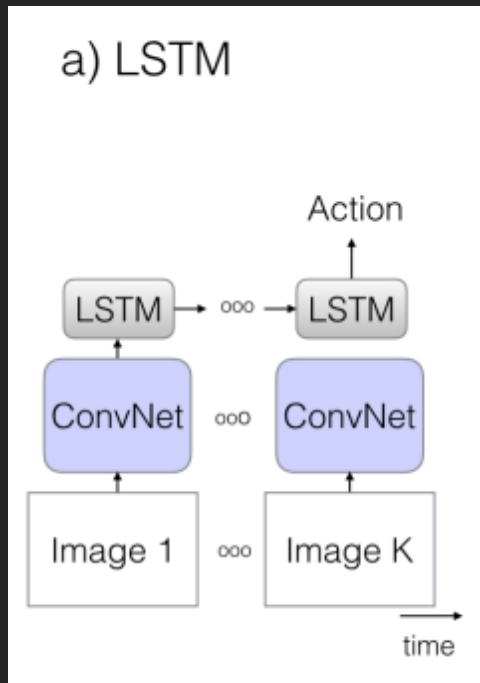
$H \times W \times C \times T$

# 3D convolutional Networks

- 3D filters at the first layer.



# Summary of Video Recognition Networks



# Pose Estimation in Videos



[https://github.com/ZheC/Realtime\\_Multi-Person\\_Pose\\_Estimation](https://github.com/ZheC/Realtime_Multi-Person_Pose_Estimation)

Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. CVPR'17

# Pose Estimation in Videos

[Demo Video:](#)

<https://www.youtube.com/watch?v=pW6nZXeWlGM&t=77s>

[https://github.com/ZheC/Realtime\\_Multi-Person\\_Pose\\_Estimation](https://github.com/ZheC/Realtime_Multi-Person_Pose_Estimation)

Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. CVPR'17

Some of my latest work:

## Temporal Relational Reasoning in Videos

Bolei Zhou, Alex Andonian, Antonio Torralba

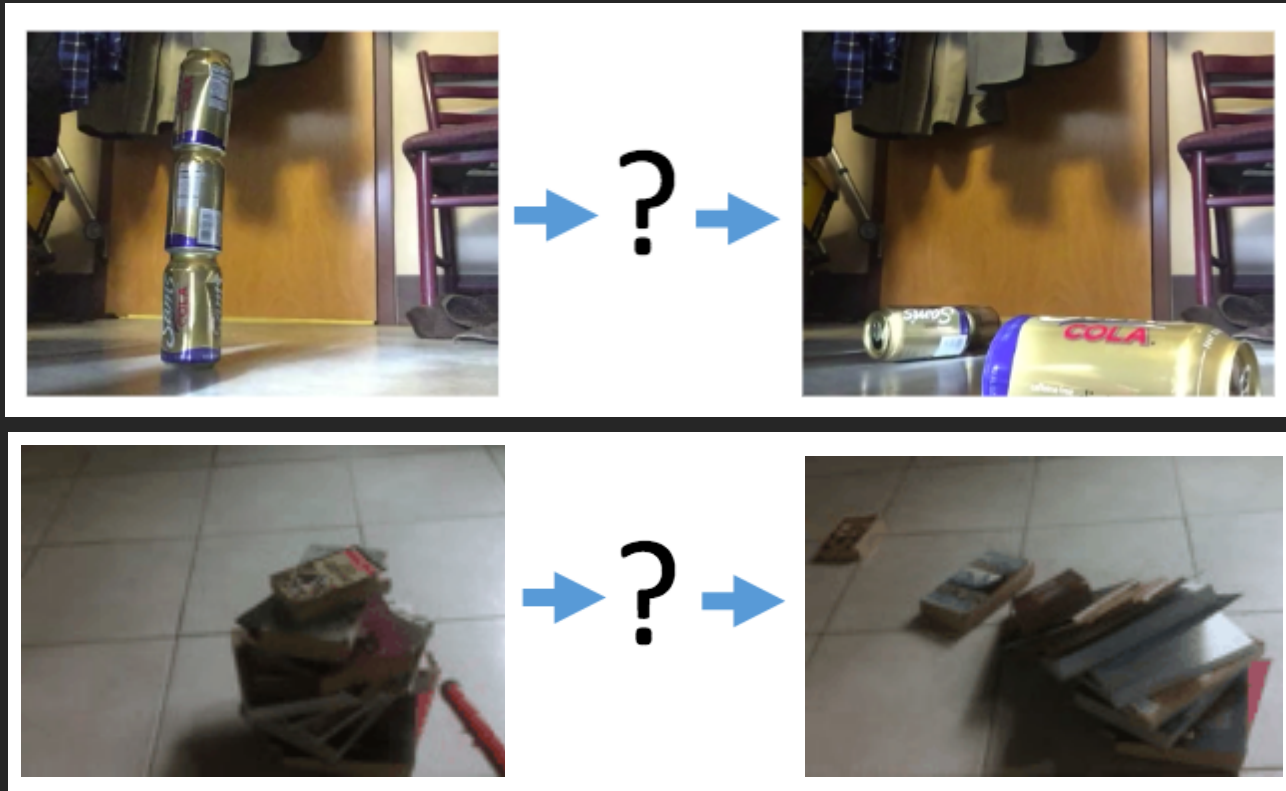
CVPR'18 submission

<https://arxiv.org/pdf/1711.08496.pdf>

# Temporal Relational Reasoning

- Infer the temporal relation between frames.

Poking a stack of something so it collapses

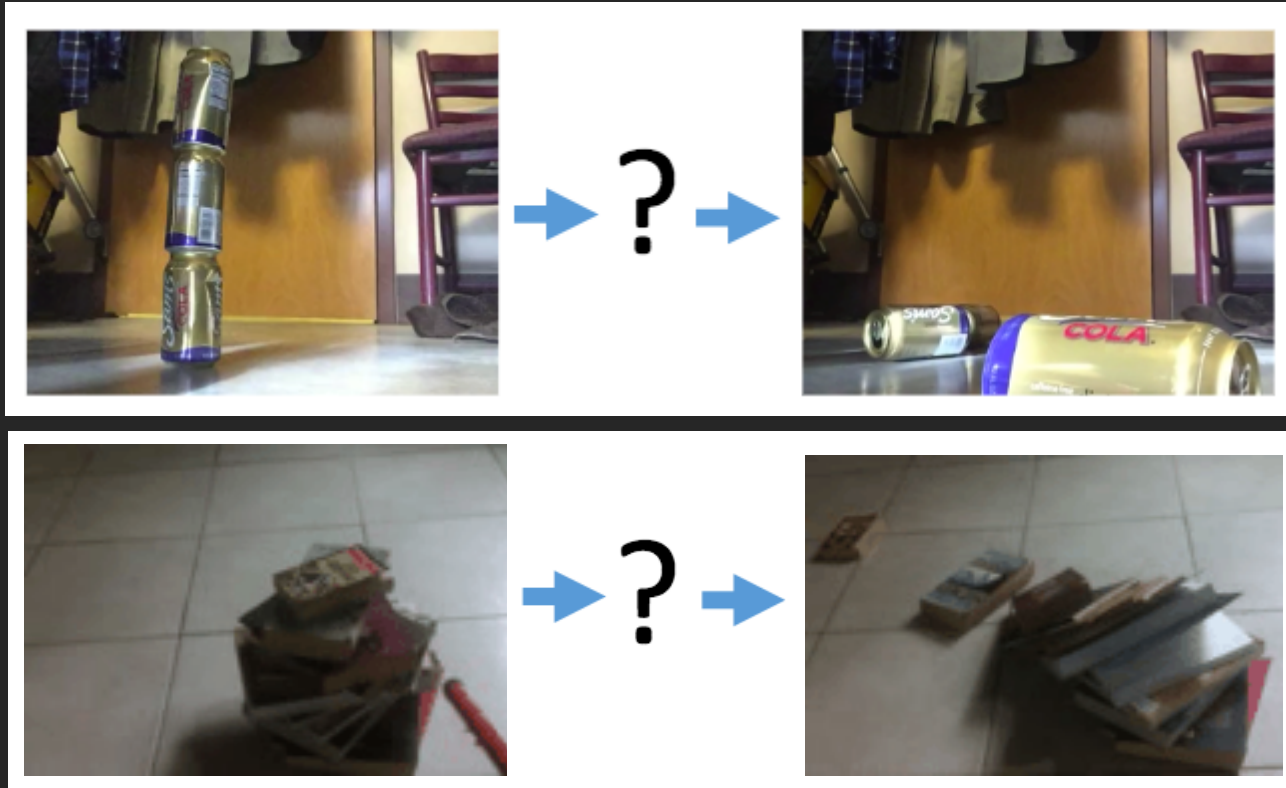




# Temporal Relational Reasoning

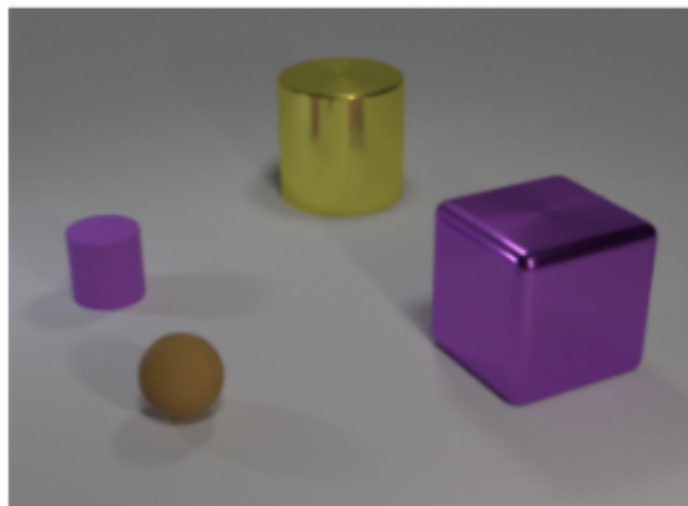
- It is the **temporal transformation/relation** that defines the activity, rather than the **appearance of objects**.

Poking a stack of something so it collapses



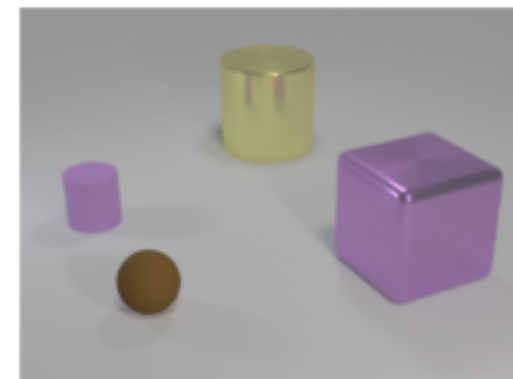
# Relational Reasoning for Visual Question Answering

**Original Image:**



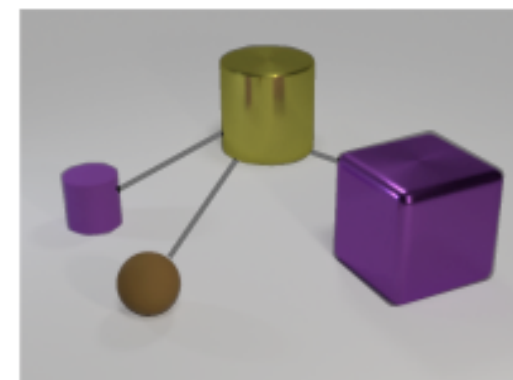
**Non-relational question:**

What is the size of the brown sphere?

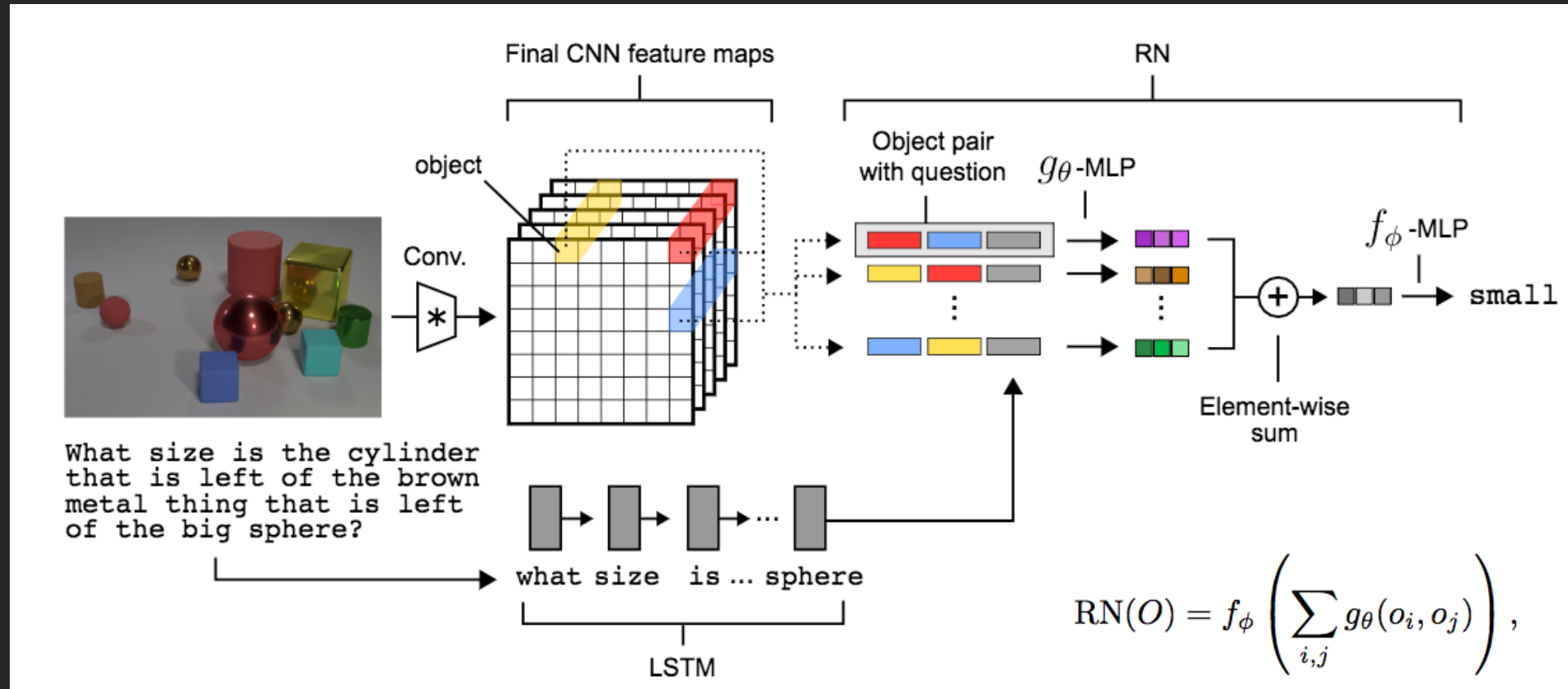


**Relational question:**

Are there any rubber things that have the same size as the yellow metallic cylinder?

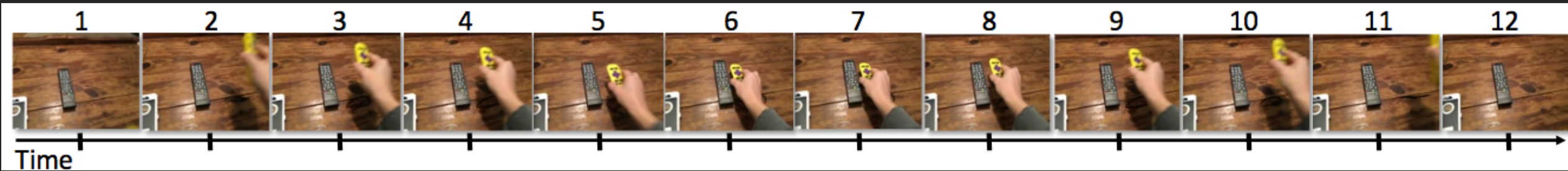


# Relational Reasoning for Visual Question Answering

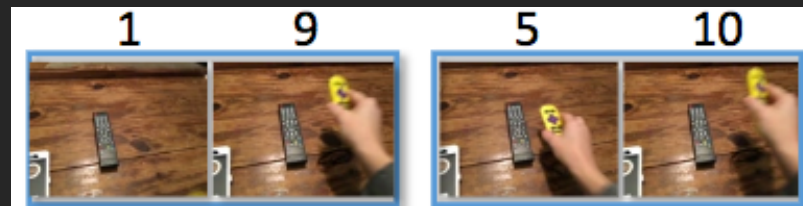


# Temporal Relations in Videos

Pretending to put something next to something



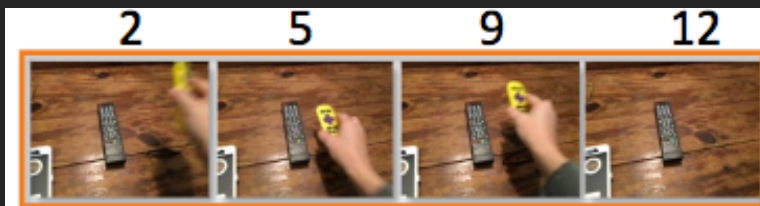
2-frame relations



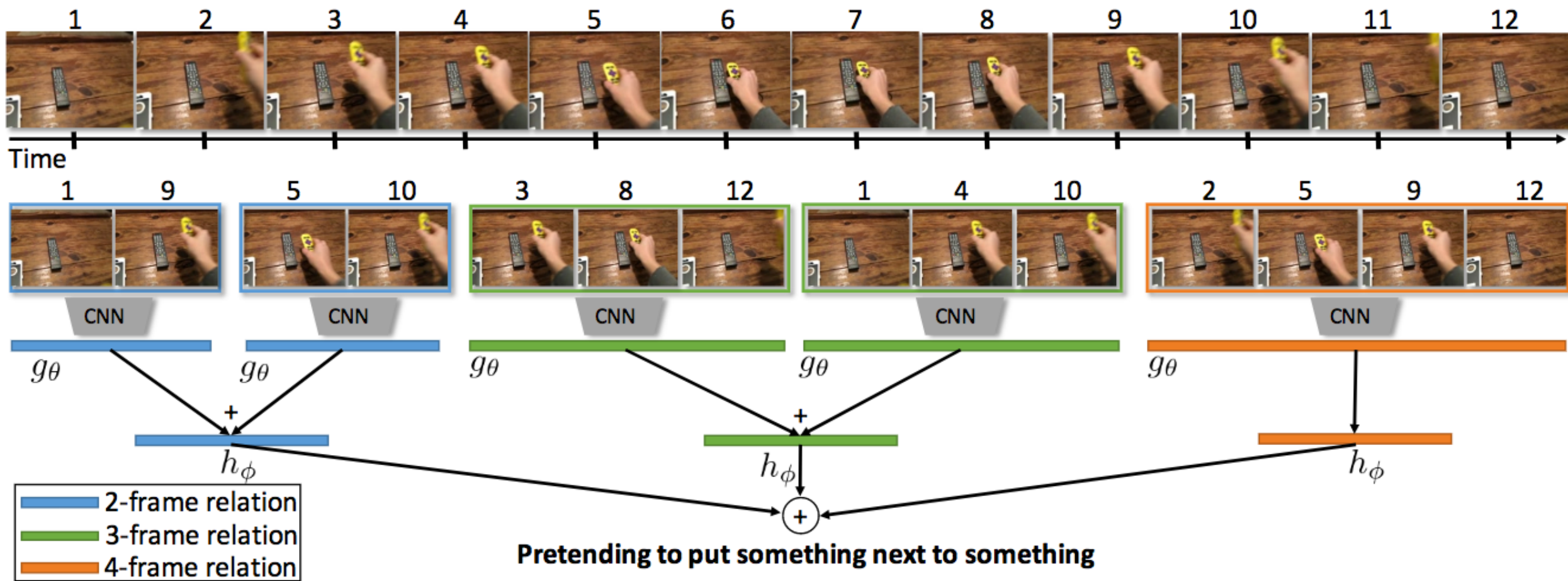
3-frame relations



4-frame relations



# Framework of Temporal Relation Networks



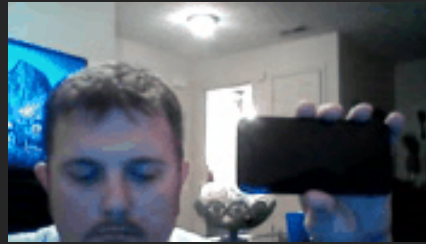
# Something-Something Dataset

- 100 K videos from 174 human-object interaction classes.

## Moving something away from something



## Plugging something into something



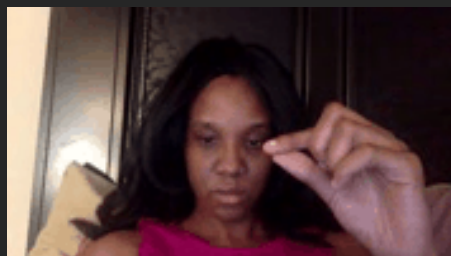
## Pulling two ends of something so that it gets stretched



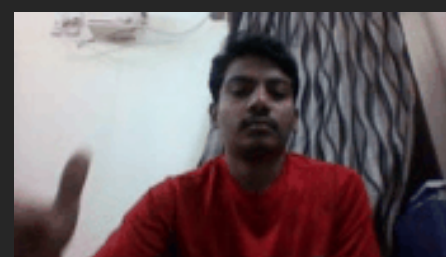
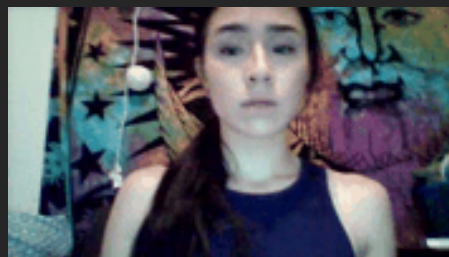
# Jester Dataset

- 140 K videos from 27 gesture classes.

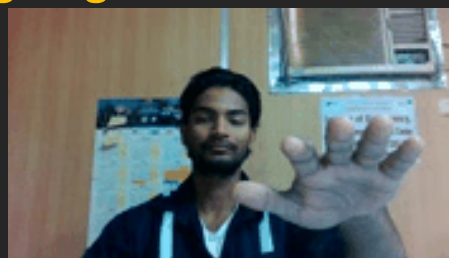
Zooming in with two fingers



Thumb down



Drumming fingers



# Experimental Results

- On Something-Something dataset

| model                    | Top1 acc.(%) | Top5 acc.    |
|--------------------------|--------------|--------------|
| single frame             | 11.41        | 33.39        |
| 2-frame TRN              | 22.23        | 48.80        |
| 3-frame TRN              | 26.22        | 54.15        |
| 4-frame TRN              | 29.83        | 58.21        |
| 5-frame TRN              | 30.39        | 58.29        |
| 7-frame TRN              | 31.01        | 59.24        |
| MultiScale TRN           | 33.01        | 61.27        |
| MultiScale TRN (10-crop) | <b>34.44</b> | <b>63.20</b> |

| model                        | Top1 acc.(%) |
|------------------------------|--------------|
| Yana Hasson                  | 25.55        |
| Harrison.AI                  | 26.38        |
| I3D by [8]                   | 27.23        |
| Guillaume Berger             | 30.48        |
| Besnet (Top1 on leaderboard) | 31.66        |
| MultiScale TRN               | <b>33.60</b> |



# Experimental Results

- On Jester dataset

| model                    | Top1 acc.(%) | Top5 acc.    |
|--------------------------|--------------|--------------|
| single frame             | 63.60        | 92.44        |
| 2-frame TRN              | 75.65        | 94.40        |
| MultiScale TRN           | 93.70        | 99.59        |
| MultiScale TRN (10-crop) | <b>95.31</b> | <b>99.86</b> |

| model                             | Top1 acc.(%) |
|-----------------------------------|--------------|
| 20BN's Jester System              | 82.34        |
| VideoLSTM                         | 85.86        |
| Guillaume Berger                  | 93.87        |
| Ford's Gesture Recognition System | 94.11        |
| Besnet (Top1 on leaderboard)      | 94.23        |
| MultiScale TRN                    | <b>94.78</b> |


# Experimental Results


- Demo Video:

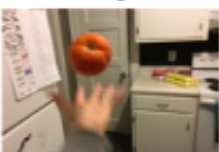
<http://relation.csail.mit.edu/>

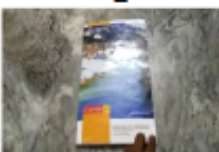
# Common sense knowledge learned by models

**Single-frame**

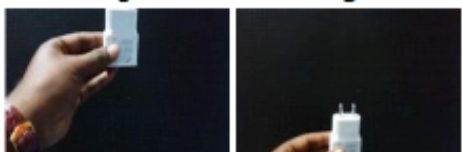
5  
Moving sth down  
  
Moving sth down (0.352)


6  
Covering sth with sth  
  
Uncovering sth (0.226)

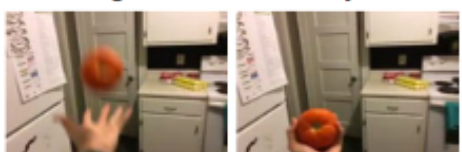
6  
Throwing sth in the air and catching it  
  
Throwing sth in the air and catching it (0.520)


2  
Pretending to open sth without actually opening it  
  
Unfolding sth (0.045)

**2-frame TRN**

0 5  
Moving sth down (0.998)  


3 6  
Covering sth (0.997)  


6 7  
Throwing sth in the air and catching it (0.999)  


2 3  
Unfolding sth (0.164)  


**3-frame TRN**

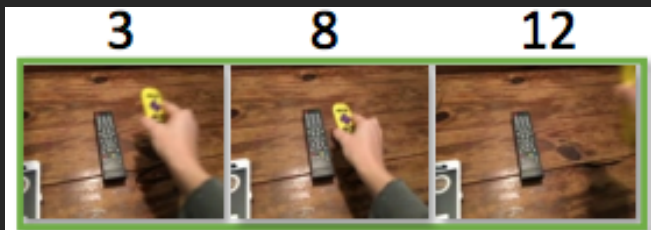
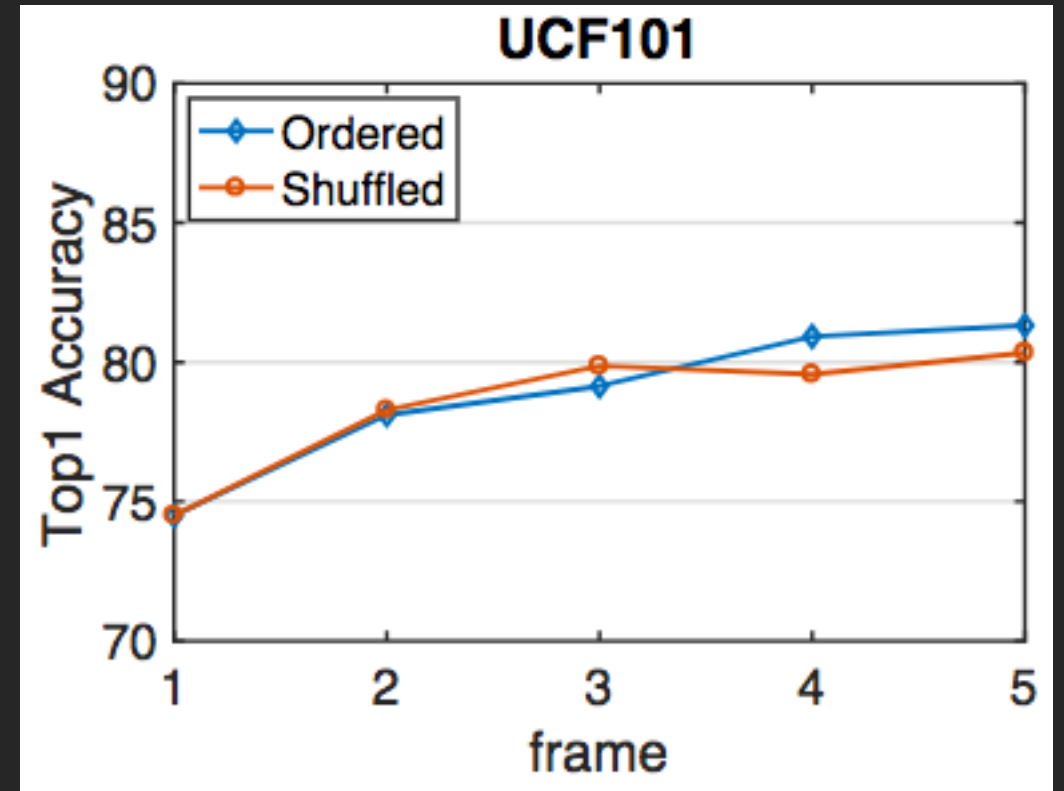
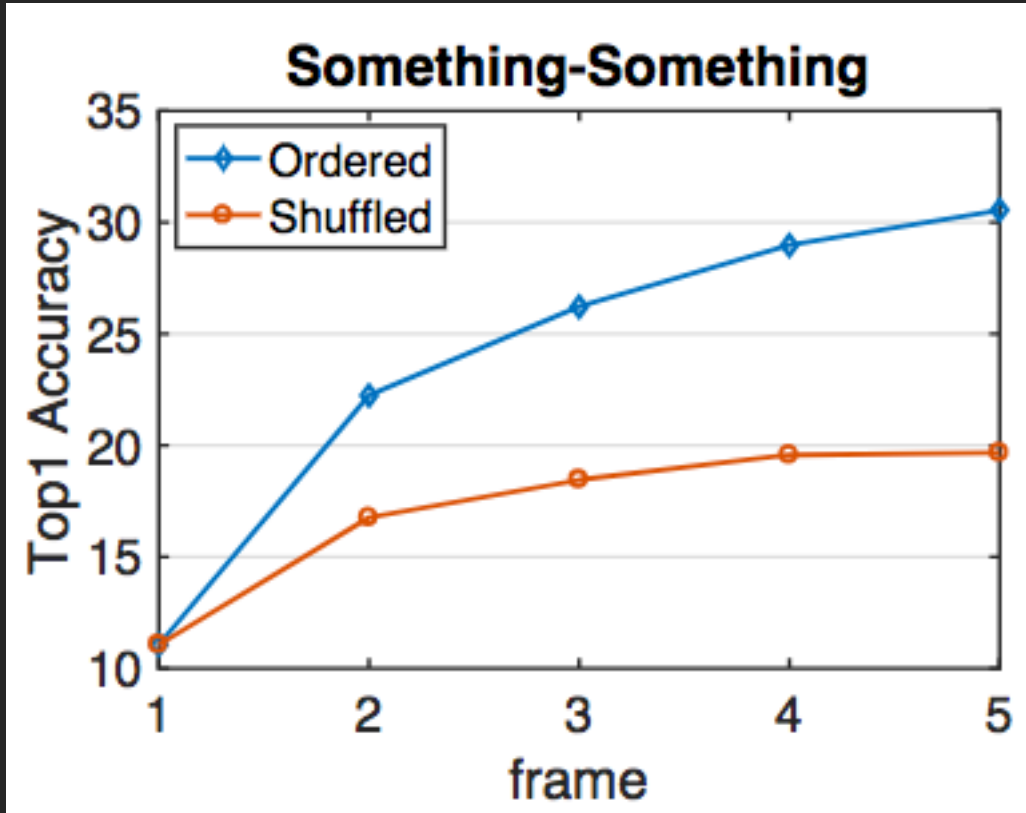
1 3 5  
Moving sth down (0.999)  


1 2 4  
Covering sth (0.998)  







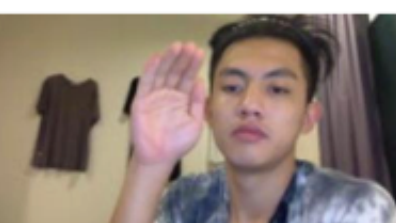
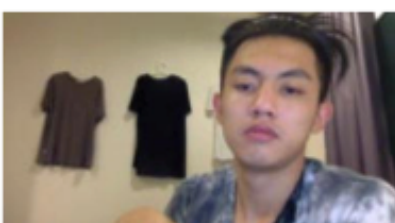

0 4 7  
Throwing sth in the air and catching it (0.999)  


2 3 6  
Pretending to open sth without actually opening it (0.828)  


# Importance of temporal orders



# Activity Forecasting

| First Frames  | Forecasts   | Ground Truth  |
|---|---|---|
|    | <ul style="list-style-type: none"><li><b>1: Tearing sth just a little bit (0.998)</b></li><li>2: Tearing sth into two pieces (0.001)</li><li>3: Pretending to be tearing sth that is not tearable (0.001)</li></ul>   |    |
|    | <ul style="list-style-type: none"><li>1: Lifting a surface with sth on it but not enough for it to slide down (0.490)</li><li><b>2: Lifting sth with sth on it (0.423)</b></li><li>3: Tilting sth with sth on it slightly so it doesn't fall down (0.079)</li></ul> |    |
|   | <ul style="list-style-type: none"><li>1: Poking sth so lightly that it doesn't or almost doesn't move (0.466)</li><li><b>2: Poking a stack of sth so the stack collapses (0.207)</b></li><li>3: Poking sth so it slightly moves (0.164)</li></ul>                   |   |
|  | <ul style="list-style-type: none"><li><b>1: Swiping Down (0.881)</b></li><li>2: Swiping Up (0.105)</li><li>3: Stop Sign (0.008)</li></ul>   |  |

# Activity Forecasting

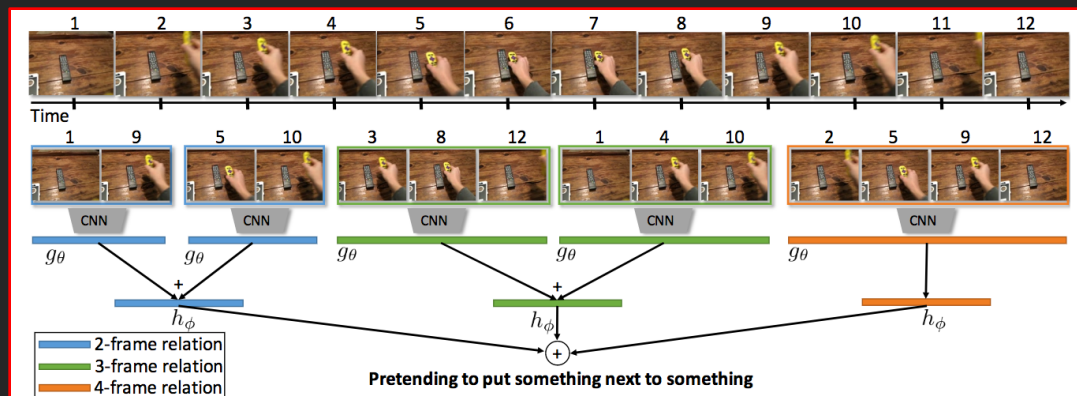
| Data      | Something |       | Jester   |       |
|-----------|-----------|-------|----------|-------|
|           | baseline  | TRN   | baseline | TRN   |
| first 25% | 9.08      | 11.14 | 27.25    | 34.23 |
| first 50% | 10.10     | 19.10 | 41.43    | 78.42 |
| full      | 11.41     | 33.01 | 63.60    | 93.70 |

# Future Directions in Activity Recognition

How to better model temporal relation?

How to make model more efficient?

- Remove the dependency on optic flow.
- Sampling of discrete frames



## Non-local Neural Networks

Xiaolong Wang<sup>1,2\*</sup>

Ross Girshick<sup>2</sup>

Abhinav Gupta<sup>1</sup>

Kaiming He<sup>2</sup>

<sup>1</sup>Carnegie Mellon University

<sup>2</sup>Facebook AI Research

### Abstract

Evolutional and recurrent operations are building blocks for processing one local neighborhood at a time. In this paper, we present non-local operations as a generic building blocks for capturing long-range dependencies inspired by the classical non-local means method.



<https://arxiv.org/pdf/1711.07971.pdf>

# Future Directions in Activity Recognition

## Activity forecasting

What's next?

kiss, hug, highfive



## Understanding long videos

Such as movie and TV shows?

